

TWO UNIVERSALITY PROPERTIES ASSOCIATED WITH THE MONKEY MODEL OF ZIPF'S LAW

RICHARD PERLINE AND RON PERLINE

ABSTRACT. The distribution of word probabilities in the monkey model of Zipf's law is associated with two universality properties: (1) the power law exponent converges strongly to -1 as the alphabet size increases and the letter probabilities are specified as the spacings from a random division of the unit interval for any distribution with a bounded density function on $[0, 1]$; and (2), on a logarithmic scale the version of the model with a finite word length cutoff and unequal letter probabilities is approximately normally distributed in the part of the distribution away from the tails. The first property is proved using a remarkably general limit theorem for the logarithm of sample spacings from Shao and Hahn, and the second property follows from Anscombe's central limit theorem for a random number of i.i.d. random variables. The finite word length model leads to a hybrid Zipf-lognormal mixture distribution closely related to work in other areas.

1. INTRODUCTION

In his popular expository article on universality examples in mathematics and physics Tao [1] mentions Zipf's empirical power law of word frequencies as lacking any "convincing explanation for how the law comes about and why it is universal." By *universality* he means the idea where systems follow certain macroscopic laws that are largely independent of their microscopic details. In this article we drop down many levels from natural language to discuss *two* universality properties associated with the toy monkey model of Zipf's law. Both of these properties were considered in Perline [2], but the presentation there was incompletely developed, and here we clarify and greatly expand upon these two ideas by showing: (1) how the model displays a nearly universal tendency towards a -1 exponent in its power law behavior; and (2) the significance of the central limit theorem (CLT) for a random number of i.i.d. variables in the model with a finite word length cutoff. We will sometimes refer to the case with a finite word length cutoff as *Monkey Twitter*, as explained in Section 3.

Our paper is laid out as follows. In Section 2 we prove the strong tendency towards an approximate -1 exponent under very broad conditions by means of a limit theorem for the logarithms of random spacing due to Shao and Hahn [3]. A somewhat longer approach to this proof is in Perline [4], where an elementary derivation of power law behavior of the monkey model is first given before applying the Shao-Hahn limit result. In Section 3, we explain the underlying lognormal structure of the *central part* of the distribution of word probabilities in the case of a *finite word length* cutoff and show how it leads to a hybrid Zipf-lognormal distribution (called a lognormal-Pareto distribution in [2]). These are universality properties in exactly Tao's sense because the distribution of the word probabilities (the macro behavior) is within very broad bounds independent of the details of how the letter probabilities for the monkey's typewriter are selected (the micro behavior). In Section 4, we discuss the

multiple connections between these results and well-known research from other areas where Pareto-Zipf type distributions have been investigated. In the remainder of this section we sketch some historical background.

It has been known for many years that the monkey-at-the-typewriter scheme for generating random “words” conforms to an inverse power law of word frequencies and therefore mimicks the inverse power form of Zipf’s [5] statistical regularity for many natural languages. *However, Zipf’s empirical word frequency law not only follows a power distribution, but as he emphasized, it also very frequently exhibits an exponent in the neighborhood of -1 .* That is, letting $f_1 \geq f_2 \geq \dots f_r \dots \geq \dots$ represent the observed ranked word frequencies, Zipf found $f_r \approx Cr^{-1}$, with $C > 0$ a constant depending on the sample size. Many qualifications have been raised regarding this approximation (including the issue of the divergence of the harmonic series); yet as I.J. Good [6] commented: “The Zipf law is unreliable but it is often a good enough approximation to demand an explanation.” Figure 1 illustrates the fascinating universal character of this word frequency law using the text from four different authors writing in four different European languages in four different centuries. The plots are shown in the log-log form that Zipf employed, and they resemble countless other examples that researchers have studied since his time. The raw text files for the four books used in the graphs were downloaded from the Public Gutenberg databank [7].

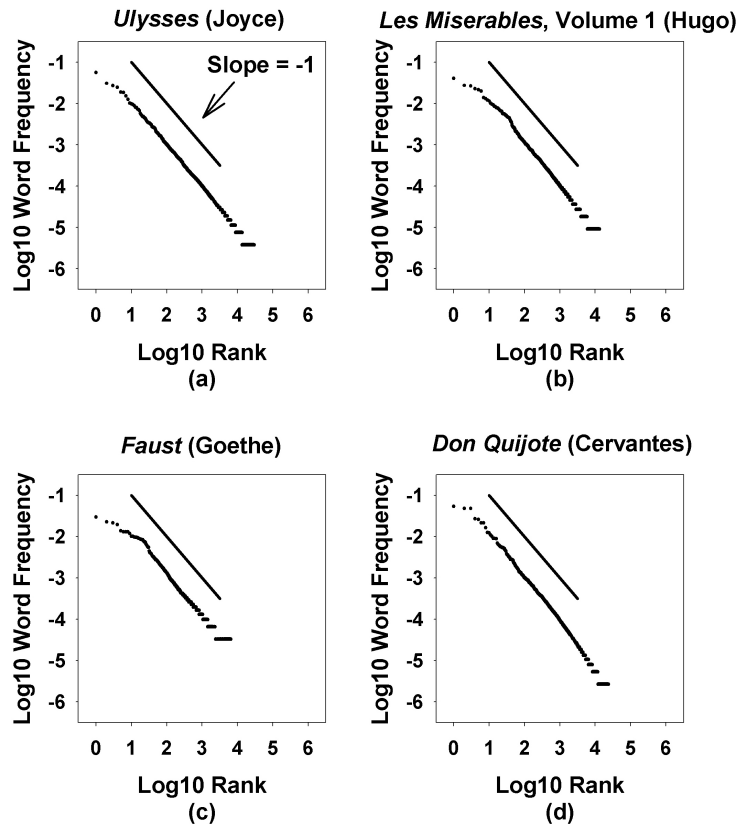


FIGURE 1. Log-log plots of relative word frequencies by rank for four authors writing in four different European languages in four different centuries. The approximate -1 slope in all the graphs is an iconic feature of Zipf’s word frequency law.

The monkey model has a convoluted history. The model can be thought of as actually embedded within the combinatoric logic of Mandelbrot's early work providing an information-theoretic explanation of Zipf's law. For example, in Mandelbrot [8] there is an appendix entitled "Combinatorial Derivation of the Law $p_r = Pr^{-B}$ " (his notation) that effectively specifies a monkey model, including a Markov version, but is not explicitly stated as such. However, his derivation there is informal, and the first completely rigorous analysis of the general monkey model with independently combined letters was given only surprisingly recently by Conrad and Michenmacher [9]. Their analysis utilizes analytic number theory and is, as a result, somewhat complicated - a point they note at the end of their article. A simpler analysis using only elementary methods based on the Pascal pyramid has now been given by Bochkarev and Lerner [10]. They have also analyzed the more general Markov problem [11] and hidden Markov models [12]. Edwards, Foxall and Perkins [13] have provided a directly relevant analysis in the context of scaling properties for paths on graphs explaining how the Markov variation can generate both a power law or a weaker scaling law, depending on the nature of the transition matrix. Perline [4] also gives a simple approach to the demonstration of power law behavior.

A clarifying and important milestone in the history of this topic came from the cognitive psychologist, G. A. Miller [14, 15]. Miller used the simple model with equal letter probabilities and an independence assumption to show that this version generates a distribution of word probabilities that conforms to an approximate inverse power law such that the probability P_r of the r^{th} largest word probability is (in a sense that needs to be made precise) of the form $P_r \approx Cr^{-\beta}$, for $-\beta < -1$. *In addition*, he made the very interesting observation that by using a keyboard of $K = 26$ letters and one space character, with the space character having a probability of .18 (about what is seen in empirical studies of English text), the value of $-\beta$ in this case is approximately -1.06, very close to the nearly -1 value in samples of natural language text as illustrated in Figure 1. In his simplified model, it turns out that $-\beta = -1 + \ln(1-s)/\ln K$, where $0 < s < 1$ is the space probability, so that $-\beta$ approaches -1 from below as K increases. (Natural logarithms are denoted $\ln(x)$ and logarithms with any other radix will be explicitly indicated, as in $\log_K(x)$. Note that Miller's β involves a *ratio* of logarithms, so that it is invariant with respect to the radix used in the numerator and denominator.) Consequently, Miller's model not only mimicks the form of Zipf's law for real languages, but with a large enough alphabet size it even mimicks the key parameter value. *In this article we give a broad generalization of Miller's observation to the case of unequal letter probabilities.* In light of our present results given in Section 2, the application of sample spacings from a uniform distribution in [2] to study the -1 behavior using an asymptotic regression line should now be viewed as just a step in a direction that ultimately led us to the Shao and Hahn limit used here.

2. PROOF THAT $-\beta$ TENDS TOWARDS -1 FROM THE ASYMPTOTICS OF LOG-SPACINGS

For our analysis we specify a keyboard with an alphabet of $K \geq 2$ letters $\{L_1, \dots, L_K\}$ and a space character S . The letter characters have non-zero probabilities q_1, q_2, \dots, q_K . The space character has probability s , so that $\sum_{i=1}^K q_i + s = 1$. A word is defined as any sequence of non-space letters terminating with a space. A word W of exactly n letters is a string such as $W = L_{i_1}L_{i_2} \dots L_{i_n}S$ and has a probability of the form $P(W) = P = q_{i_1}q_{i_2} \dots q_{i_n}s$ because letters are struck independently. The space character with no preceding letter character will be considered a word of length zero. The rank ordered sequence of descending word

probabilities in the ensemble of all possible words is written $P_1 = s > P_2 \geq \dots P_r \geq \dots$ ($P_1 = s$ is always the first and largest word probability.) We break ties for words with equal probabilities by alphabetical ordering, so that each word probability has a unique rank r .

Conrad and Mitzenmacher [9] give a carefully constructed definition of power law behavior in the monkey model as the situation where there exist two positive constants C_1 and C_2 such that the inequality

$$C_1 r^{-\beta} \leq P_r \leq C_2 r^{-\beta} \quad (\beta > 0) \quad (1)$$

holds for sufficiently large r . Following the argument in Bochkarev and Lerner [10], β in the monkey model turns out to be simply the solution to

$$q_1^{1/\beta} + q_2^{1/\beta} + \dots + q_K^{1/\beta} = 1, \quad (2)$$

with $\beta > 1$.¹ In the Miller model with equal letter probabilities, $q_1 = \dots = q_K = (1-s)/K$, so from equation (2), β in this case is found to be $1 - \ln(1-s)/\ln K$. In the Fibonacci example given by Conrad and Mitzenmacher [9] and in Mitzenmacher [16] they use $K = 2$ letters with probabilities $q_1, q_2 = q_1^2$ and $q_1 < (-1 + \sqrt{5})/2$ so that $q_1 + q_2 < 1$. Then β is the solution to $q_1^{1/\beta} + q_1^{2/\beta} = 1$, giving

$$\beta = \frac{\ln q_1}{\ln((-1 + \sqrt{5})/2)}. \quad (3)$$

To understand conditions leading to $-\beta \approx -1$, we define spacings through a random division of the unit interval and then state the Shao-Hahn limit law. Let X_1, X_2, \dots, X_{K-1} be a sample of $K-1$ i.i.d. random variables drawn from a distribution on $[0, 1]$ with a bounded density function $h(x) \leq M < \infty$.² Write the order statistics of the sample as $X_{1:K-1} \geq X_{2:K-1} \geq \dots \geq X_{K-1:K-1}$. The K spacings D_i are defined as the differences between the successive order statistics: $D_1 = 1 - X_{1:K-1}$, $D_i = X_{i-1:K-1} - X_{i:K-1}$ for $2 \leq i \leq K-1$ and $D_K = X_{K-1:K-1}$. We'll refer to this as a *generalized broken stick process*. By Shao and Hahn [3] Corollary 3.6, we have

$$\frac{1}{K} \sum_{i=1}^K \ln(K D_i) \xrightarrow{a.s.} - \int_0^1 h(x) \ln h(x) dx - \lambda, \quad (4)$$

as $K \rightarrow \infty$ and where *a.s.* signifies *almost sure convergence*, $\lambda = .577\dots$ is the Euler constant and $-\int_0^1 h(x) \ln h(x) dx$ is the differential entropy of $h(x)$. Clearly,

$$\frac{1}{K} \sum_{i=1}^K \ln(K D_i) = \ln K + \frac{1}{K} \sum_{i=1}^K \ln D_i, \quad (5)$$

so dividing through by $\ln K$ gives

¹Bochkarev and Lerner [10] indicate $\ln p(r) - \ln r/\gamma$ in the inequality for their Theorem 1, but what was evidently intended is $\ln p(r) + \ln r/\gamma$. Their $1/\gamma$ is equal to our β .

²Shao and Hahn present the conditions for this limit in a more general way that reduces to this simpler statement when a density function exists.

$$\frac{\ln K}{\ln K} + \frac{1}{K} \frac{\sum_{i=1}^K \ln D_i}{\ln K} \xrightarrow{a.s.} \frac{-\int_0^1 h(x) \ln h(x) dx}{\ln K} - \frac{\lambda}{\ln K}, \quad (6)$$

as $K \rightarrow \infty$. The right side of the limit in (6) goes to 0 because $-\int_0^1 h(x) \ln h(x) dx / \ln K$ goes to 0 by the boundedness of the density $h(x)$. Expressing logarithms with a radix $= K$ then leads to the limit

$$\frac{\sum_{i=1}^K \log_K D_i}{K} \xrightarrow{a.s.} -1 \text{ as } K \rightarrow \infty. \quad (7)$$

Our *universality property* for β will now follow almost immediately from this. We use sample spacings to populate the K letter probabilities for the monkey keyboard. Since s is the probability of the space character, define $q_i = (1-s)D_i$ ($1 \leq i \leq K$) so that $\sum_{i=1}^K q_i = 1-s$. Let $\bar{m}_K = \sum_{i=1}^K \log_K q_i / K$. Then from the limit (7), we have

$$\begin{aligned} \bar{m}_K &= \sum_{i=1}^K \frac{\log_K q_i}{K} \\ &= \log_K(1-s) + \frac{\sum_{i=1}^K \log_K D_i}{K} \xrightarrow{a.s.} -1 \end{aligned} \quad (8)$$

as $K \rightarrow \infty$. Since $\bar{m}_K \xrightarrow{a.s.} -1$ and $-\beta < -1$, showing that $\bar{m}_K \leq -\beta$ will prove that $-\beta \xrightarrow{a.s.} -1$ as $K \rightarrow \infty$. To see that this is the case, note:

$$\sum_{i=1}^K \frac{\log_K q_i^{1/\beta}}{K} = \log_K \left[(q_1^{1/\beta} q_2^{1/\beta} \cdots q_K^{1/\beta})^{1/K} \right] \quad (9)$$

$$\leq \log_K \left(\frac{q_1^{1/\beta} + q_2^{1/\beta} \cdots + q_K^{1/\beta}}{K} \right) \quad (10)$$

$$= \log_K \frac{1}{K} = -1, \quad (11)$$

where the inequality in (10) follows from the geometric-arithmetic mean inequality. Therefore, from

$$\sum_{i=1}^K \frac{\log_K q_i^{1/\beta}}{K} = \frac{1}{\beta} \bar{m}_K \leq -1, \quad (12)$$

we see that $\bar{m}_K \leq -\beta$ and so $-\beta \xrightarrow{a.s.} -1$. In the special case of Miller's model using all equal letter probabilities, $\bar{m}_K = -\beta$.

The broad generality of the Shao-Hahn limit leads to the near-universal characterization of the approximation $-\beta \approx -1$ as K increases. Figure 2 presents graphical results using simulations of sample spacings from several distributions to illustrate this phenomenon. In Figure 2(a) we plot $\log_{10} P_r$ by $\log_{10} r$ for the Miller model with exactly the parameter values he used, i.e., $K = 26$ letters, a space probability $s = .18$ and equal letter probabilities $q_1 = \cdots = q_{26} = (1 - .18)/26$. The graph is based on the ranks $1 \leq r \leq 475255 = \sum_{i=0}^4 26^i$, corresponding to all words of length ≤ 4 non-space letters. Figures 2(b) through 2(d) are

based on generating $K = 26$ letter probabilities from three different continuous distributions with bounded densities $h(x)$ on $[0, 1]$ using a generalized broken stick method to obtain the sample spacings. Again $s = .18$ was used for the probability of the space character and the letter probabilities were populated with values from the spacings with $q_i = (1 - s)D_i$, $1 \leq i \leq 26$, so that in each case, $\sum_{i=1}^{26} q_i = .82$, as in Miller's example. For each graph in Figure 2(b) - 2(d), we generated the largest 475255 word probability values in order to match the Miller example of Fig 2(a). The three continuous distributions, all defined on $[0, 1]$, are:

- a uniform distribution with density $h(x) = 1$;
- a beta $B(3, 2)$ distribution with density

$$h(x) = \frac{\Gamma(3+2)}{\Gamma(3)\Gamma(2)} x^{3-1}(1-x)^{2-1}; \quad (13)$$

- a triangular distribution with density

$$h(x) = \begin{cases} 4x, & \text{if } 0 \leq x \leq 1/2 \\ 4(1-x) & \text{if } 1/2 \leq x \leq 1. \end{cases} \quad (14)$$

The graphs in Figures 2(b)-2(d) illustrate our theoretical derivation, but we also note that the very linear plots indicate what appears to be an almost “immediate” convergence to power law behavior exceeding what might be expected from the asymptotics of Conrad-Mitzenmacher and Bochkarov-Lerner.

3. MONKEY TWITTER: ANSCOMBE'S CLT FOR THE MODEL WITH A FINITE WORD LENGTH CUTOFF

The discussion of the finite word length version of the monkey model in [2] was incomplete because it did not explicitly provide the normalizing constants for the application of Anscombe's CLT for a random number of i.i.d. variables, as we do here. We also want to explain more clearly the nature of the hybrid Zipf-lognormal distribution that results from a finite length cutoff when letter probabilities are not identical.

The focus in Section 2 on the infinite ensemble of word probabilities \mathbf{P}_∞ ranked to give $P_1 > P_2 \geq P_3 \dots$ has actually obscured the significant *hierarchical structure* of the monkey model - what Mandelbrot [17] (p. 345), called a *lexicographic tree* - which only becomes evident when word length is considered. To see this, write $\mathbf{P}_{\text{len}=n}$ for the multiset of all word probabilities for monkey words of exactly length n . There are K^n probabilities in $\mathbf{P}_{\text{len}=n}$ having a total sum of $(q_1 + q_2 + \dots + q_K)^n s = (1 - s)^n s$ and an average value of $((1 - s)/K)^n s$, declining geometrically with n . Now taking word length into account, Mandelbrot's lexicographic tree structure becomes clear with the simple case of $K = 2$ letters: the root node of the tree has a probability of s (for the “space” word of length 0); it has two branches to the next level consisting of nodes for words of length 1 with probabilities $q_1 s$ and $q_2 s$; these each branch out to the next level with four nodes having probability values $q_1^2 s, q_1 q_2 s, q_2 q_1 s, q_2^2 s$ and so on.

The relevance of the CLT is seen first by representing any probability $P \in \mathbf{P}_{\text{len}=n}$ as a product of i.i.d. random variables times the constant s : $P = X_1 X_2 \dots X_n s$, where each X_i takes on one of the letter probability values q_1, q_2, \dots, q_K with probability $1/K$ (i.e., we use the natural counting measure to construct a probability space on $\mathbf{P}_{\text{len}=n}$). Let $\mu_1 = \sum_{i=1}^K \ln q_i / K$ and $\sigma_1^2 = \sum_{i=1}^K (\ln q_i - \mu_1)^2 / K$. Assume that $\sigma_1^2 > 0$, i.e., the letter probabilities are not all equal. Then $\mu_1 + \ln s$ is the mean and σ_1^2 is the variance of $\ln P$ for $P \in \mathbf{P}_{\text{len}=1}$, and it follows

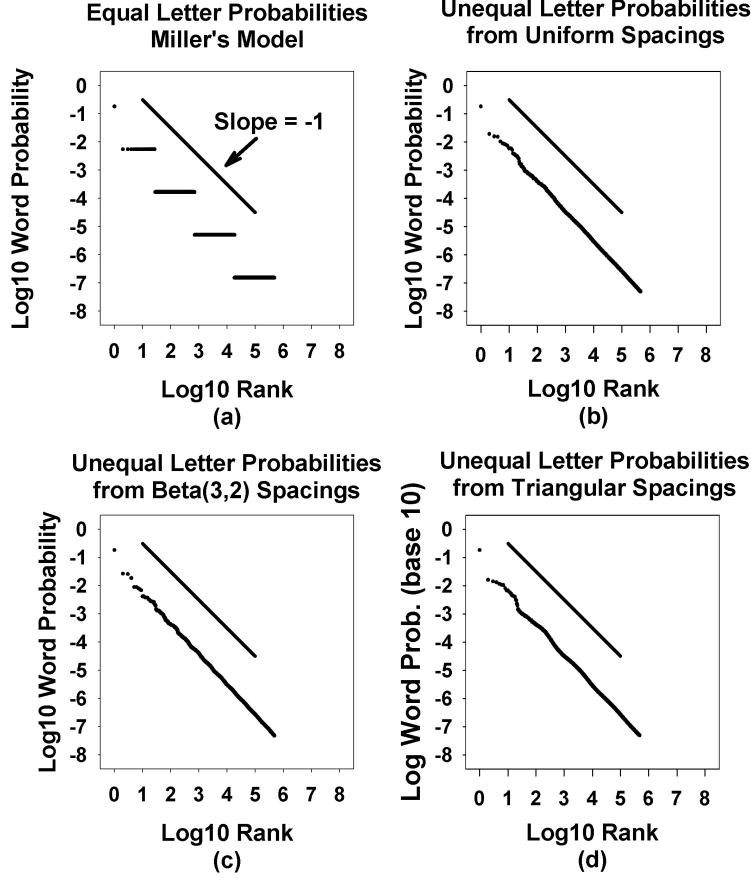


FIGURE 2. Log-log plots of monkey word probabilities by rank showing the asymptotic tendency towards a -1 exponent (slope on the log-log scales) using four different distributions to generate letter probabilities: equal probabilities in 2(a) and a generalized broken stick process in 2(b)-2(d). $K = 26$ letters were used in all cases and the largest $\sum_{i=0}^4 26^i = 475255$ word probabilities are displayed.

that $n\mu_1 + \ln s$ and $n\sigma_1^2$ are the respective mean and variance of $\ln P = S_n = \sum_{i=1}^n \ln X_i + \ln s$ for $P \in \mathbf{P}_{\text{len}=n}$. Therefore, $(S_n - n\mu_1)/\sqrt{n\sigma_1^2}$ is asymptotically normally distributed $N(0, 1)$ (the term $\ln s$ is rendered negligible in the asymptotics) so that for sufficiently large n , P itself will be approximately lognormal $LN(n\mu_1, n\sigma_1^2)$. This approximate normality for the log-probabilities of words of fixed length is quite obvious and was noted in passing by Mandelbrot [18] (p. 210). However, he missed a *stronger* observation that the $N_n = \sum_{i=0}^n K^i$ word probabilities for $P \in \mathbf{P}_{\text{len} \leq n} = \cup_{i=0}^n \mathbf{P}_{\text{len}=i}$ have a distribution that behaves in its central part away from the tails very much like $P \in \mathbf{P}_{\text{len}=n}$. That is, a version of the CLT can be applied to words of length $\leq n$, not just to words of length n , and the two distributions, *in one sense*, are very close to each other.

To explain why this is so, for $P \in \mathbf{P}_{\text{len} \leq n}$ represent $\ln P = S_{R_n} = \sum_{i=1}^{R_n} \ln X_i + \ln s$ as a sum of a *random number* of i.i.d. random variables (plus the constant $\ln s$), where R_n is itself a random variable with a finite geometric distribution, $\text{Prob}\{R_n = i\} = K^i/N_n$, ($1 \leq i \leq n$). We will not repeat the calculations of [2] here, but it is easily shown that $R_n/N_n \rightarrow 1$ in

probability. Because this limiting constant is 1, Anscombe’s generalization of the CLT [19] (Theorem 3.1, p. 15) can be applied with the *identical* normalizing constants $n\mu_1$ and $\sqrt{n\sigma_1^2}$ as used above with $\ln P$ for $P \in \mathbf{P}_{\text{len}=\mathbf{n}}$. Consequently, it is *also* true that for $P \in \mathbf{P}_{\text{len}\leq\mathbf{n}}$, the normalized sum $(S_{R_n} - n\mu_1)/\sqrt{n\sigma_1^2}$ has an asymptotic $N(0, 1)$ distribution. In other words the two random sums S_n and S_{R_n} behave so similarly in their centers that the word probabilities in $\mathbf{P}_{\text{len}=\mathbf{n}}$ and $\mathbf{P}_{\text{len}\leq\mathbf{n}}$ are *both* approximately $LN(n\mu_1, n\sigma_1^2)$. However, it is essential to remark that the two distributions have very different upper tail behavior. In this regard, Le Cam’s [20] comment that French mathematicians use the term “central” referring to the CLT “because it refers to the center of the distribution as opposed to its tails” is particularly relevant.

The behavior of $P \in \mathbf{P}_{\text{len}\leq\mathbf{n}}$ is illustrated in the graphs of Figure 3(a) and 3(b). The plot in Figure 3(a) was generated just as in Figure 2(b) except that a finite length cutoff of $n \leq 4$ letters has been applied. To make this clear, in Figure 2(b) the plot shows the top 475255 word probabilities in \mathbf{P}_∞ generated using $K = 26$ letters derived from uniform spacings. In contrast, using the same letter probabilities, the plot in Figure 3(a) shows the 475255 word probabilities generated with word length ≤ 4 letters, i.e., all the values of $P \in \mathbf{P}_{\text{len}\leq 4}$. Word length in the case of *unequal* letter probabilities is certainly correlated, *but not perfectly*, with word probability: words of shorter length will, *on average*, have a higher probability than those of longer length, but except in Miller’s degenerate case, there will always be reversals where a longer word will have a higher probability than a shorter word.³ However, writing $P_{1:N_n} > P_{2:N_n} \cdots \geq P_{N_n:N_n}$ for the ranked values of $\mathbf{P}_{\text{len}\leq\mathbf{n}}$, it should be evident that for any given rank r , $P_{r:N_n} \leq P_r$ and that $P_{r:N_n} = P_r$ when n is sufficiently large. In short, $\mathbf{P}_{\text{len}\leq\mathbf{n}}$ inherits its upper tail power law behavior from \mathbf{P}_∞ , which is illustrated by comparing the top part of the curve in Figure 3(a) with the corresponding part of Figure 2(b).

Figure 3(b) uses the same data points from Figure 3(a) graphed as a normal quantile plot, and its roughly linear appearance for the logarithm of word probabilities for $P \in \mathbf{P}_{\text{len}\leq\mathbf{n}}$ conforms to an approximate Gaussian fit, although certainly the bending in the upper half of the distribution departs a bit from the linear trend of the lower half. The fact that a distribution can have a power law tail and lognormal central part and *still look lognormal over essentially its entire range* may seem surprising. The discussion in [21] about *power law mimicry* in the upper tail of lognormal distributions helps to explain why this can happen. We will call the distribution for $\mathbf{P}_{\text{len}\leq\mathbf{n}}$ a Zipf-lognormal distribution.

The monkey model with a fixed word length cutoff can prove useful as a motivating idea. In the next section, we will discuss how models with the same branching tree structure have been proposed many times in the past, typically in settings where something like a finite word length is natural to consider. For the moment, think of the social networking service, Twitter, which allows members to exchange messages limited to at most 140 characters. Define *Monkey Twitter* with a finite limit of $n + 1$ characters. For convenience, in any implementation of a Monkey Twitter random experiment we will assume that monkeys would always fill up their allotted message space of $n + 1$ characters. Monkey words still require a terminating space character, so it is possible for a monkey to type (at most) one “non-word,” which can vary in length from 1 to $n + 1$ characters, and will always be the last part of a message string. Non-words are discarded, and the probabilities for the legitimate

³In natural languages, Zipf [5] underscored that “the length of a word tends to bear an inverse relationship to its relative frequency,” which he called the *Law of Abbreviation*. In many respects, this is the starting point for his *Principle of Least Effort*.

monkey words, varying in length from 0 to n non-space characters plus a terminating space character, will correspond to the values in $P_{\text{len} \leq n}$.

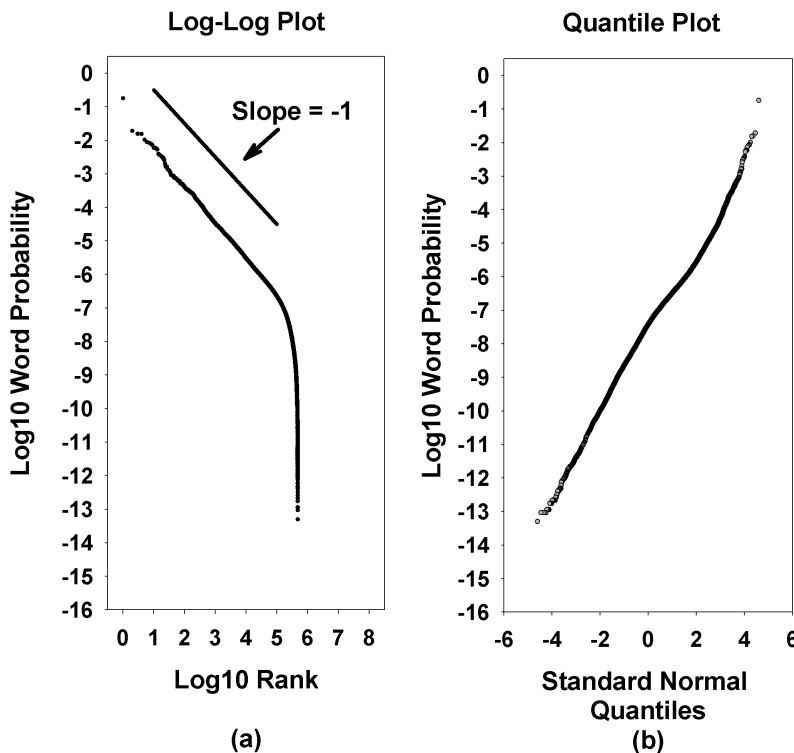


FIGURE 3. Figure 3(a) is a log-log plot of monkey word probabilities by rank for all words of length ≤ 4 non-space characters using letter probabilities from uniform spacings. The linear upper tail coincides closely with the previous Figure 2(b), but the power law clearly breaks down. Figure 3(b) shows the same word probabilities in a standard normal quantile plot. Its rough linearity confirms an approximate Gaussian fit over the whole distribution even though the upper tail is a power law as seen on the left in Figure 3(a). We refer to this distribution as a Zipf-lognormal hybrid, and it has many connections to other distributions discussed in the statistical literature.

4. CONNECTIONS TO OTHER WORK

The subject of power law distributions is a vast topic of research reaching across diverse scientific disciplines [16, 21, 22, 23]. Here we provide a brief sketch of how our results connect to some other work and how they can be considered in a much more general light.

The monkey model viewed in terms of its hierarchical tree structure is the starting point for understanding its general nature. Though Miller did not explicitly present his model as a branching tree, several researchers at almost the same time in the mid-to-late 1950s were highlighting this form using the equivalent of equal letter probabilities to motivate the occurrence of empirical power law distribution in other areas. For example, Beckman [24] introduced this idea in connection with the power law of city populations within a country

(“Auerbach’s law” [25]). Using essentially the same logic as Miller, but in a completely different setting, Beckman assumed that a given community will have K satellite communities, each with a constant decreasing fraction of the population at a higher level. That is, if A is the population of the largest city and $0 < p < 1$, there will be K nearby satellite communities with population pA , K^2 smaller communities nearby to these with population p^2A , and so on. Mandelbrot [26] (p. 226) had an apt expression for this pattern: he described it as “compensation between two exponentials” because the number of observations at each level increases geometrically while the mean value decreases geometrically. Beckman then went on to note that if instead of using a constant decreasing fraction p , one used a random variable X on $(0, 1)$, the population at the n^{th} level down would be a random variable of the form $X_1X_2 \cdots X_nA$, leading to an approximate lognormal distribution of populations at this level. This corresponds to our discussion of monkey word probabilities in $\mathbf{P}_{\text{len}=n}$, although Beckman was not aware of the still stronger statement of lognormality for the probabilities of words of length $\leq n$ in $\mathbf{P}_{\text{len} \leq n}$ that we demonstrated using Anscombe’s CLT.

Many more examples of a branching tree structure essentially equivalent to Miller’s monkey model with equiprobable letters have been proposed over the years to motivate the occurrence of a huge variety of empirical power law distributions, including such size distributions as lake areas, island areas (“Korcak’s law”), river lengths, etc. Indeed, Mandelbrot’s [17] classic book on fractals is a rich source of these. This equiprobability case is so simple to analyze that it has been invoked over and over again to illustrate how power law behavior can arise. However, the more complicated case using unequal letter probabilities (or proportions) *and a finite word length cutoff* is what really underscores the close analogy between the monkey Zipf-lognormal distribution and other mixture models exhibiting hybrid power law tails and approximately lognormal central ranges.

Empirical distributions of this form have appeared with great frequency. In fact there is a long history of researchers, including Pareto, first discovering what appears to be an empirical power law over an entire distribution because they start out by looking at only the largest observations (cities, corporations, islands, etc.). However, when they extend their measurements to the part of the distribution below the upper tail, *which is always more difficult to observe*, the power law behavior typically breaks down. Several examples of this are chronicled in [21]. The power law for city sizes noted above is a perfect illustration. Auerbach in 1913 [25] looked at the 94 largest German cities from a 1910 census and showed a good power law fit, but because he did not include the thousands of smaller communities, he may not have been aware that this relationship breaks down. Modern studies of the *full* distribution of communities, such as Eeckhout’s [27] analysis of the populations of 25,359 places from U.S. Census 2000, prove with high confidence that the *bulk* of the community populations fit an approximate lognormal distribution and that the power law behavior is confined to the upper tail.

Montroll and Shlesinger [28, 29] explained how to generate a Pareto-lognormal hybrid for income distributions by using a hierarchical mixture model of lognormal distributions. This is motivated from several angles, including the notion that higher classes amplify their income by organizing in such a way as to benefit from the efforts of others. Reed and Hughes [30] have provided a far-ranging framework for understanding these hybrid distributions across the entire spectrum of disciplines where they have been discovered: “physics, biology, geography, economics, insurance, lexicography, internet ecology, etc.” In [30] he and Hughes give a condensed summary showing that “if stochastic processes with exponential

growth in expectation are killed (or observed) randomly, the distribution of the killed or observed state exhibits power-law behavior in one or both tails.” This work encompasses: (1) geometric Brownian motion (GBM); (2) discrete multiplicative processes; (3) homogeneous birth and death processes; and (4) Galton-Watson branching processes. In one of many variations on this theme, in his GBM model [31] Reed specifies a stochastic process that uses lognormal distributions varying continuously in time with an exponential mixing distribution and shows that this generates a “Double Pareto-Lognormal Distribution.” This has an asymptotic Pareto upper tail and a central part approximately lognormal; in addition, it exhibits interesting asymptotic behavior in the it lower tail where it is characterized by a direct (rather than an inverse) power law. Reed has gone to great effort to demonstrate the high quality of the fit of this distribution for all ranges of values (low, middle, high) across a wide variety of size distributions such as particle and oil field sizes [31], incomes [32], internet file sizes and the sizes of biological genera [30].

In today’s nomenclature, the term *Zipf’s law* has come to mean any power law distribution with an exponent close to the value -1 , not just the word frequency law. (To be clear here, when we refer to power law distributions, we also mean the hybrids with power law tails that we have been considering.) The subset of power laws with this restricted exponent value is surprisingly large and includes the distribution of firm sizes [33], city sizes [34], the famous Gutenberg-Richter law for earthquake magnitudes [35] and many others.

Gabaix [34] pointed out that while it has long been known that stochastic growth processes could generate power laws, “these studies stopped short of explaining why the Zipf exponent should be 1.” To address this question, he has given a theoretical explanation of the genesis of a -1 exponent for the Zipf law for city populations, but in fact, his approach can be regarded more generally. His key idea is a variation of Gibrat’s [36] *Law of Proportion Effect*: using a fixed number of normalized city population quantities that sum to 1 and assuming growth processes (expressed as percentages) randomly distributed as i.i.d. variables with a common mean and variance, he proves that a steady state will satisfy Zipf’s law. This proof requires an additional strong assumption in order to reach a steady state, namely, a lower bound for the (normalized) population size. Gabaix goes on to discuss relaxed versions of his model and how it fits into the larger context of similar work done by others.

In their monograph on Zipf’s law Saichev et al [37] modify and extend the core Gabaix idea in numerous ways that render it more realistic. For concreteness, they present their work using the terminology of financial markets and corporate asset values, but they are very clear on how their models are relevant to a broad range of “physical, biological, sociological and other applications.” As with Reed, GBM plays an important role in their investigations, which incorporate birth and death processes, acquisitions and mergers, the subtleties of finite-size effects and other features. Notably, they focus on the sets of conditions that lead to an approximate -1 exponent.

In both Harr  moes and Tops  e [38] and Corominas-Murtra and Sol   [39] -1 emerges as a consequence of certain information theoretic ideas pertaining to the entropy function under growth assumptions. In [38] the focus is on language and the expansion of vocabulary size while simultaneously maintaining finite entropy. The perspective in [39] is broader, but still based on the idea of systems with an expanding number of possible states “characterized by special features on the behavior of the entropy.”

There is another topic area of statistical research that impacts on our work. Natural language word frequency distributions have been characterized as *Large Number of Rare*

Events (LNRE) distributions because when collections of text are examined, no matter how big, there are always a large number of words that occur very infrequently in the sample [40]. LNRE behavior indicates a tip-of-the-iceberg phenomenon resulting from a bias that captures the most common words, but necessarily misses a vast quantity of rarely used words. This is a classic and much studied problem encountered in species abundance surveys in ecological research, where a key question becomes estimating the size of the *zero abundance class* - the typically great number of species not observed in a sample because of their rarity [41].

To get an idea of the significance of this issue, consider that Figure 3(a) graphs the *population* (or *parent* or *theoretical*) word probabilities for the Zipf-lognormal distribution, and not the *sample* frequencies as would be obtained from actually carrying out the Monkey Twitter experiment. Simple visual inspection of Figure 3(a) indicates that the linear part of the graph (i.e., the power law tail) holds for about 5 logarithmic decades or about the first 100000 word probabilities out of the total of 475255 probabilities plotted there. However, these 100000 probabilities comprise 97.4% of the total probability mass of all 475255 values. Intuitively, it should be clear that sampling from this Zipf-lognormal population distribution will be dominated by the Zipf part, not the lognormal part. Unless the sample size was astronomically large, so that the large number of low probability words showed up, the underlying structure of the parent distribution would not reveal itself. To carry this idea still further, imagine the situation if the word length cutoff was on the order of $n = 140$ characters, such as with real Twitter. No experiment could be run within any realistic time frame to ever hope to obtain a sample sufficiently large to uncover the true hybrid character of the population distribution - the sample would always appear as a Zipf law, not a Zipf-lognormal law. This kind of *visibility bias* has been a constant and recurring theme in all areas where Pareto-Zipf type distributions have been studied [2, 21]. We believe its significance has been poorly appreciated in relation to this topic.

Finally, we will take a more bird's-eye view of matters and remark that our application of random spacings is very much in the spirit of the enormously fruitful study of *random matrix theory* and a class of stochastic systems referred to as the *KPZ universality class*. Along these lines Borodin and Gorin [42] have discussed a variety of probabilistic systems "that can be analyzed by essentially algebraic methods," yet are applicable to a broad array of topics. Miller's demonstration of power law behavior and an approximate -1 exponent with increasing K for the monkey model is in a similar vein. We regard this result as analogous to the Borodin and Gorin example of the De Moivre-Laplace proof of the CLT for a sequence of i.i.d. Bernoulli trials, which depends on having an explicit pre-limit distribution and then taking "the limit of the resulting expression." (We want to add, however, that Miller's model is ultimately *too simple* to reveal the approximate lognormal behavior of the monkey word probabilities - for that, we needed to assume non-identical letter probabilities.) Our point is to make clear that the monkey model fits into a much larger conceptual scheme than appears at first glance.

5. CONCLUSIONS

Sample spacings provide a natural way to populate the letter probabilities in the monkey model through a random division of the unit interval. The Shao and Hahn asymptotic limit law for the logarithms of spacings then leads to the result that the exponent in the power law of word frequencies will tend towards a -1 value under broad conditions as the number of letters in the alphabet increases. The monkey model can also be viewed as a

branching tree structure. In that light, Anscombe's CLT reveals an underlying lognormal central part of the word frequency distribution when a finite word length cutoff is imposed. The resulting hybrid Zipf-lognormal distribution has many connections to other work. The visibility bias inherent in sampling from this distribution is similar to what has been noted as a characteristic of many different empirical power laws.

REFERENCES

- [1] Tao, T. E pluribus unum: from complexity, universality. In *The Best Writing on Mathematics 2013*; Pitici, M., Ed.; Princeton University Press: Princeton, NJ, USA, 2014; pp. 32-46.
- [2] Perline, R. Zipf's law, the central limit theorem and the random division of the unit interval. *Phys. Rev. E* **1996**, *54*, 220-223.
- [3] Shao, Y.; Hahn, M.G. Limit theorems for the logarithm of sample spacings. *Stat. Prob. Lett.* **1995**, *24*, 121-132.
- [4] Perline, R. The random division of the unit interval and the approximate -1 exponent in the monkey-at-the-typewriter model of Zipf's law. Submitted to *Stat. Prob. Lett.* March 5, 2015 and under review.
- [5] Zipf, G.K. *Human Behavior and the Principle of Least Effort*; Addison-Wesley: Cambridge, MA, USA, 1949.
- [6] Good, I.J. Statistics of language: introduction. In *Encyclopedia of Linguistics, Information and Control*; Meetham, A. R., Ed.; Pergamon Press: NY, NY, USA, 1969; pp. 567-581.
- [7] Project Gutenberg. Michael Hart. Available online: <http://www.gutenberg.org/> (accessed on June 6, 2014). The plain text files for these four books, each in their original language, were used to obtain the word frequencies plotted in Figure 1: Joyce, J. *Ulysses*; Hugo, V. *Les Misérables*, Volume 1; Goethe, W. *Faust*; Cervantes, M. *Don Quijote*.
- [8] Mandelbrot, B.B. On recurrent noise limiting coding. In *Information Networks, the Brooklyn Polytechnic Institute Symposium*; Weber, E., Ed.; Interscience: NY, NY, USA, 1955; pp. 205-221.
- [9] Conrad, B.; Mitzenmacher, M. Power laws for monkeys typing randomly: the case of unequal letter probabilities. *IEEE Trans. Info. Theory*, **2004**, *50*, 1403-1414.
- [10] Bochkarev, V.V.; Lerner, E. Yu. The Zipf law for random texts with unequal probabilities of occurrence of letters and the Pascal pyramid. *Russian Math.*, **2012**, *56*, 25-27.
- [11] Bochkarev, V.V.; Lerner, E. Yu. Strong power and subexponential laws for an ordered list of trajectories of a Markov chain. *Electronic Journal of Linear Algebra*, **2014**, *27*, 534-556.
- [12] Bochkarev, V.V.; Lerner, E. Yu. Zipf exponent of trajectory distribution in the hidden Markov model. *Journal of Physics: Conference Series*, **2014**, *490*, 012008.
- [13] Edwards, R.; Foxall, E.; Perkins T.J. Scaling properties of paths on graphs. *Electronic Journal of Linear Algebra*, **2012**, *23*, 966-988.
- [14] Miller, G.A. Some effects of intermittent silence. *Amer. J. Psych.* **1957**, *70*, 311-314.
- [15] Miller, G.A.; Chomsky, N. Finitary Models of Language Users. In *Handbook of Mathematical Psychology*; Luce, R.D.; Bush, R.R.; Galanter, E., Eds.; Wiley: NY, NY, USA, 1963; Volume 2, pp. 419-491.
- [16] Mitzenmacher, M. A brief history of generative models for power law and lognormal distributions. *Internet Math.*, **2003**, *1*, 226-251.
- [17] Mandelbrot, B.B. *The Fractal Geometry of Nature*; W.H. Freeman and Company, NY, NY, USA, 1983.
- [18] Mandelbrot, B.B. On the theory of word frequencies and on related markovian models of discourse. In *Structure of Language and Its Mathematical Aspects: Proceedings of Symposia on Applied Mathematics Volume 3*; Jakobson, R., Ed.; Amer. Math. Soc.: Providence, RI, USA, 1961; pp. 190-219.
- [19] Gut, A. *Stopped Random Walks: Limit Theorems and Applications*; Springer-Verlag: NY, NY, USA, 1988.
- [20] Le Cam, L. The central limit theorem around 1935. *Statistical Science*, **1986**, *1*, 78-96.
- [21] Perline, R. Strong, weak and false inverse power laws. *Statistical Science*, **2005**, *20*, 68-88.
- [22] Clauset, A.; Shalizi, C.R.; Newman, M.E. Power law distributions in empirical data. *SIAM Review*, **2009**, *51*, 661-703.
- [23] Arnold, B. *Pareto Distributions - Second Edition*; CRC Press, Taylor and Francis: Boca Raton, FL, USA, 2015.

- [24] Beckman, M.J. City hierarchies and the distribution of city sizes. *Economic Development and Cultural Change*, **1958**, 6, 243-248.
- [25] Auerbach, F. Das Gesetz der Bevölkerungskonzentration. *Petermanns Geographische Mitteilungen*, **1913**, 59, 74-76.
- [26] Mandelbrot, B.B. *Fractals and Scaling in Finance: Discontinuity, Concentration, Risk SELECTA VOL-UME E*; Springer: NY, NY, USA, 1997; pp. 219-251.
- [27] Eeckhout, J. Gibrat's law for (all) cities. *American Economic Review*, **2004**, 94, 1429-1451.
- [28] Montroll, E.; Shlesinger, E. On $1/f$ noise and other distributions with long tails. *Proc. Natl Acad. Sci. USA*, **1982**, 79, 3380-3383.
- [29] Montroll, E.; Shlesinger, E. Maximum entropy formalism, fractals, scaling phenomena, and $1/f$ noise: a tale of tails. *J. Stat. Phys.*, **1983**, 32, 209-230.
- [30] Reed, W.J.; Hughes, B. D. From gene families and genera to incomes and internet file sizes: why power laws are so common in nature. *Phys. Rev. E*, **2002**, 66, 067103.
- [31] Reed, W.J. On Pareto's law and the determinants of Pareto exponents. *J. Income Dist.*, **2004** 13,1-2.
- [32] Reed, W.J. The double Pareto-lognormal distribution – a new parametric model for size distributions. *Comm. in Stat.*, **2004**, 33, 1733-1753.
- [33] Axtell, R.L. Zipf distribution of U.S. firm sizes. *Science*, **2001**, 293, 1818-1820.
- [34] Gabaix, X. Zipf's law for cities: an explanation. *The Quarterly Journal of Economics*, **1999**, 114, 739-767.
- [35] Kagan, Y.Y. Universality of the seismic moment-frequency relations. *Pure and Applied Geophysics*, **1999**, 155, 537-573.
- [36] Gibrat, R. *Les inegalites economiques*; Librairie du Recueil Sirey, Paris, France, 1931.
- [37] Saichev, A.; Malevergne, Y.; Sornette, D. *Theory of Zipf's Law and Beyond*; Springer-Verlag: Berlin, Germany, 2010.
- [38] Harr  moes, P.; Tops  oe, F. Maximum entropy fundamentals. *Entropy*, **2001**, 3, 191-226.
- [39] Corominas-Murtra, B.; Sol  , R. Universality of Zipf's law. *Physical Review E*, **2010**, 82, 011102.
- [40] Baayen, R. H. *Word Frequency Distributions*, Kluwer Academic Publishers, Dordrecht, Holland, 2001.
- [41] Bunge, J.; Fitzpatrick, M. Estimating the number of species: a review. *J. Amer. Stat. Assoc.*, **1993**, 88, 364-373.
- [42] Borodin, A.; Gorin, V. Lectures on integrable probability. Available online: <http://arxiv.org/pdf/1212.3351v2.pdf> (accessed on 1 November 2015).

FLEXIBLE LOGIC ANALYTICS, 34-50 80TH STREET, JACKSON HEIGHTS, NY 11372, USA
E-mail address: richperline@gmail.com

DEPARTMENT OF MATHEMATICS, DREXEL UNIVERSITY, PHILADELPHIA, PA, 19104, USA
E-mail address: ronald.k.perline@drexel.edu